

Supplementary Note: Assessment of TumorBoost based on
tumor/normal pair TCGA-02-0001 in the Affymetrix
GenomeWideSNP_6 data set TCGA,GBM,CRMAv2

Henrik Bengtsson, Pierre Neuvial, Terence P. Speed

March 5, 2010

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 2 |
| 2 | Data set | 2 |
| 2.1 | Preprocessing methods | 2 |
| 2.2 | Stratification on genotype confidence scores | 2 |
| 2.3 | List of change points | 2 |
| 3 | Region: TCGA-02-0001:Chr2@35-74,cp=57+/-1,s=1/2 | 3 |
| 3.1 | Decrease in Heterozygosity (DH) and total copy-number tracks | 3 |
| 3.2 | Allele B fraction density plots | 4 |
| 3.3 | ROC curves | 4 |
| 3.4 | (β_N, β_T) plots | 5 |
| 3.5 | Allele-specific copy number estimates | 6 |
| 4 | Region: TCGA-02-0001:Chr2@75-110,cp=96+/-1,s=1/4 | 7 |
| 4.1 | Decrease in Heterozygosity (DH) and total copy-number tracks | 7 |
| 4.2 | Allele B fraction density plots | 8 |
| 4.3 | ROC curves | 8 |
| 4.4 | (β_N, β_T) plots | 9 |
| 4.5 | Allele-specific copy number estimates | 10 |
| 5 | Region: TCGA-02-0001:Chr2@100-130,cp=110+/-1,s=4/0 | 11 |
| 5.1 | Decrease in Heterozygosity (DH) and total copy-number tracks | 11 |
| 5.2 | Allele B fraction density plots | 12 |
| 5.3 | ROC curves | 12 |
| 5.4 | (β_N, β_T) plots | 13 |
| 5.5 | Allele-specific copy number estimates | 14 |
| 6 | Region: TCGA-02-0001:Chr13@0-70,cp=45+/-1,s=0/2 | 15 |
| 6.1 | Decrease in Heterozygosity (DH) and total copy-number tracks | 15 |
| 6.2 | Allele B fraction density plots | 16 |
| 6.3 | ROC curves | 16 |
| 6.4 | (β_N, β_T) plots | 17 |
| 6.5 | Allele-specific copy number estimates | 18 |
| 7 | Bootstrap estimates of test statistics for all regions | 19 |
| A | Data files | 21 |
| A.1 | Total copy numbers | 21 |
| A.2 | Allele B fractions | 21 |
| A.3 | Genotype calls | 21 |
| B | Session information | 22 |

1 Introduction

This report, which is automatically generated, assesses the performance of the TumorBoost method based on a few change points in a particular tumor/normal pair. For more details on the evaluation methods, see the main TumorBoost manuscript.

2 Data set

The evaluation is this report is based on the tumor/normal pair (01C,10A) for individual TCGA-02-0001 in the data set TCGA,GBM,CRMAv2.

2.1 Preprocessing methods

The data was generated on the Affymetrix GenomeWideSNP_6 chip type.

There is one CEL file per hybridization. Each CEL file was preprocessed separately using CR-MAv2 Bengtsson *et al.* (2009), without relying neither on reference samples nor prior estimates.

2.2 Stratification on genotype confidence scores

No stratification on genotype confidence scores is done. All heterozygous SNPs are used in this evaluation.

2.3 List of change points

For this data set, we have selected a few regions for which one can safely assume that there exists a single copy number change point. By definition, each change point separates two sets of genomic loci such that the true Decrease in Heterozygosity (DH) is the same within one set of loci but differs between the two sets. These regions were selected visually. For each region we chose a large enough safety margin to make our evaluation independent of the uncertainty on the true location of the change point.

| Chr | Region | Change point | Margin | Before | After |
|-----|---------|--------------|--------|--------------|----------------|
| 2 | 35-74 | 57 | 1 | gain (1,2) | deletion (0,1) |
| 2 | 75-110 | 96 | 1 | gain (1,2) | gain (1,3) |
| 2 | 100-130 | 110 | 1 | gain (1,3) | normal (1,1) |
| 13 | 0-70 | 45 | 1 | normal (1,1) | deletion (0,1) |

Table 1: Regions in TCGA-02-0001 used for the evaluation and that each contain a single changepoint. All positions and lengths are in units of Mb.

We next compare how well each of these change points is detected using the above preprocessed signals followed or not by TumorBoost normalization using the ROC analysis described in the main TumorBoost manuscript at the full resolution as well as smoothed resolution with bin sizes $h = \{1, 2, 4\}$. Specifically, we compare the following three methods: (1) “**raw**”: preprocessed signals without TumorBoost normalization. (2) “**TBN,NGC**”: preprocessed signals followed by TumorBoost normalization using NGC genotype calls. For completeness we also include an evaluation of Total copy numbers (TCN).

3 Region: TCGA-02-0001:Chr2@35-74,cp=57+/-1,s=1/2

3.1 Decrease in Heterozygosity (DH) and total copy-number tracks

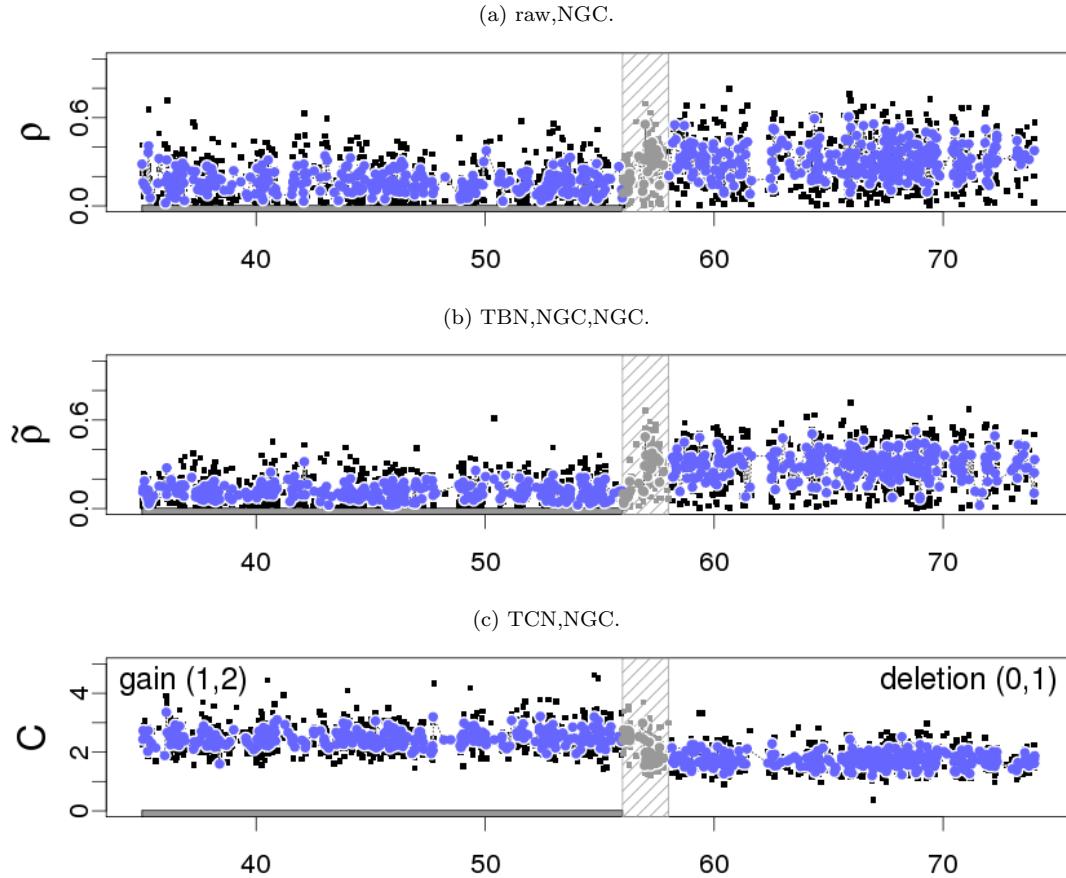


Figure 1: Decrease in Heterozygosity (DH) and total copy numbers for region TCGA-02-0001:Chr2@35-74,cp=57+/-1,s=1/2. Only heterozygous SNPs are plotted. There are 1219 loci of state gain (1,2) ("negatives") and 1219 loci of state deletion (0,1) ("positives"), where the latter are highlighted with a solid bar beneath. In total 129 loci within the safety margin were excluded.

3.2 Allele B fraction density plots

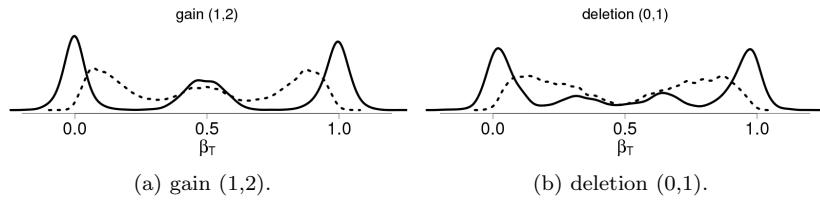


Figure 2: Density of raw (dashed lines) and TumorBoost-normalized (solid lines) allele B fractions for region TCGA-02-0001:Chr2@35-74,cp=57+/-1,s=1/2.

3.3 ROC curves

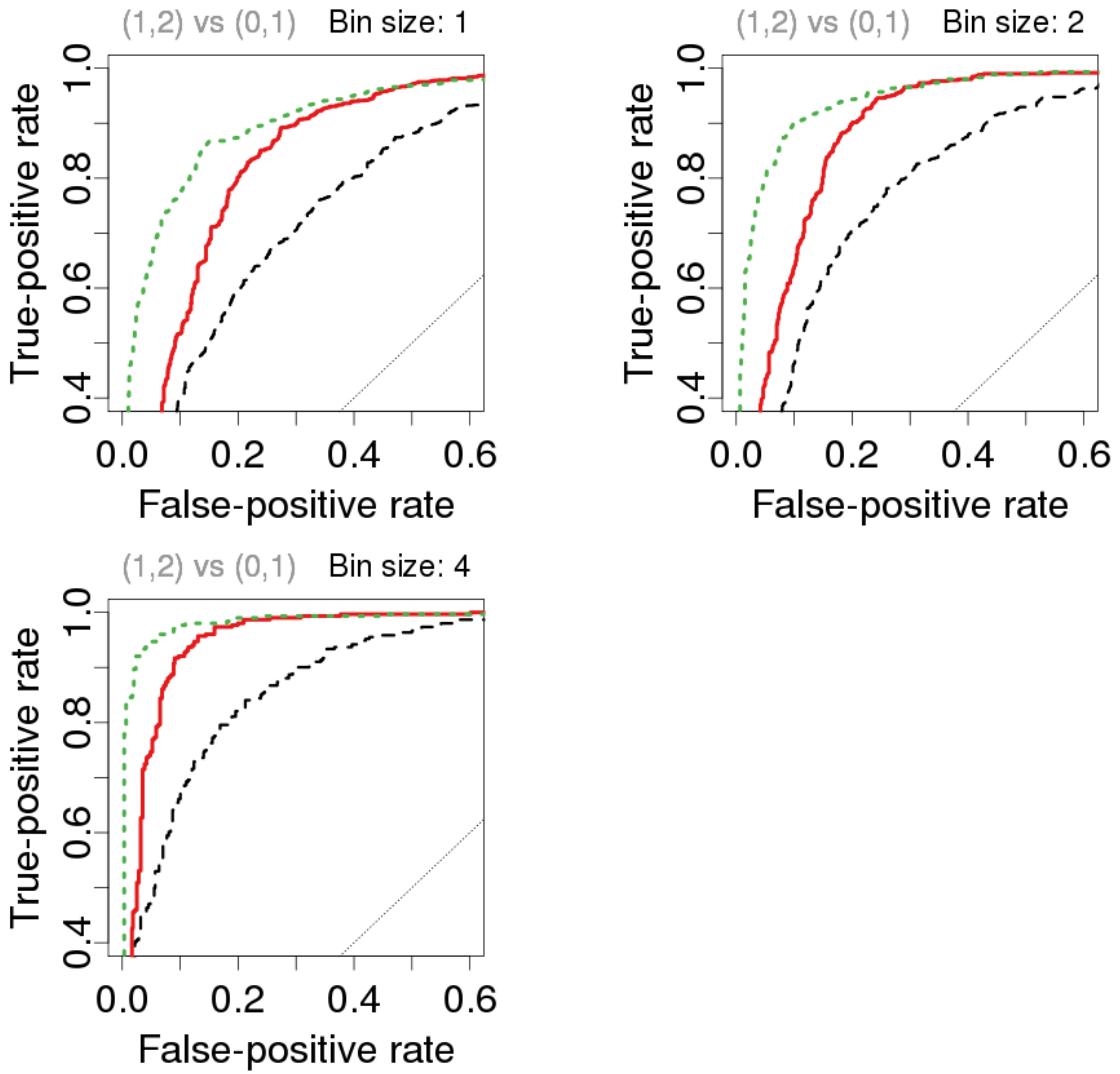


Figure 3: ROC curves for each preprocessing method at the full resolution as well as 2 different amounts of smoothing (using the `mean()` function) for region TCGA-02-0001:Chr2@35-74,cp=57+/-1,s=1/2. Legend: raw,NGC (dashed; #000000), TBN,NGC,NGC (solid; #E41A1C) and TCN,NGC (dotted; #4DAF4A).

3.4 (β_N, β_T) plots

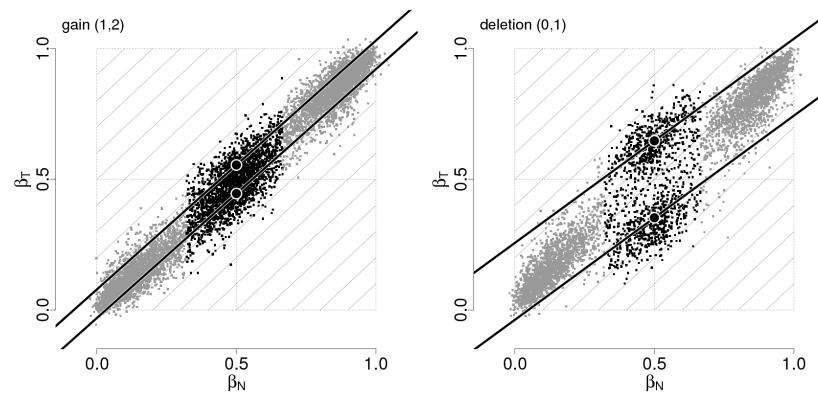


Figure 4: raw,NGC for region TCGA-02-0001:Chr2@35-74,cp=57+/-1,s=1/2.

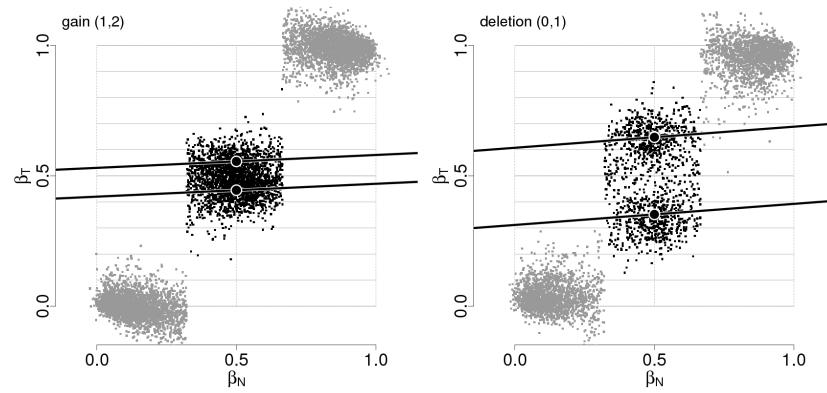


Figure 5: TBN,NGC,NGC for region TCGA-02-0001:Chr2@35-74,cp=57+/-1,s=1/2.

3.5 Allele-specific copy number estimates

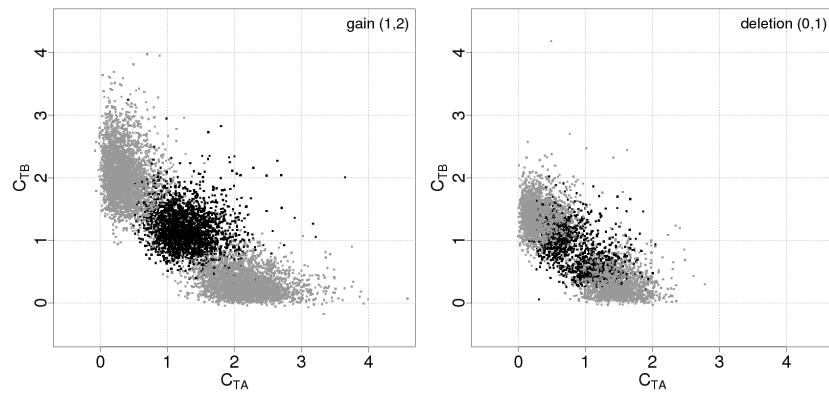


Figure 6: raw,NGC for region TCGA-02-0001:Chr2@35-74,cp=57+/-1,s=1/2.

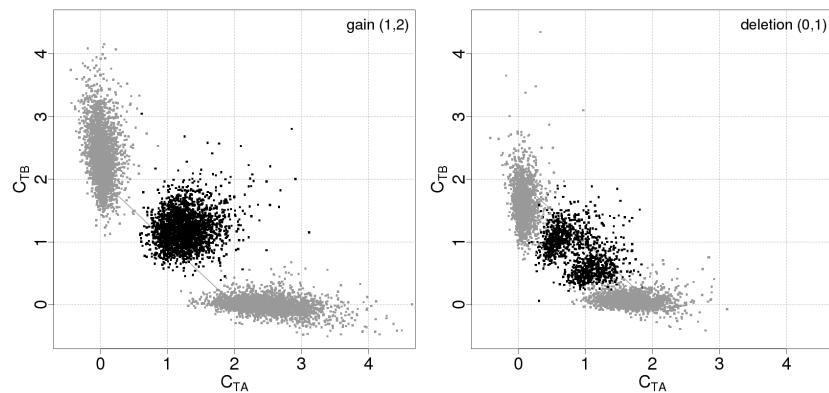


Figure 7: TBN,NGC,NGC for region TCGA-02-0001:Chr2@35-74,cp=57+/-1,s=1/2.

4 Region: TCGA-02-0001:Chr2@75-110,cp=96+/-1,s=1/4

4.1 Decrease in Heterozygosity (DH) and total copy-number tracks

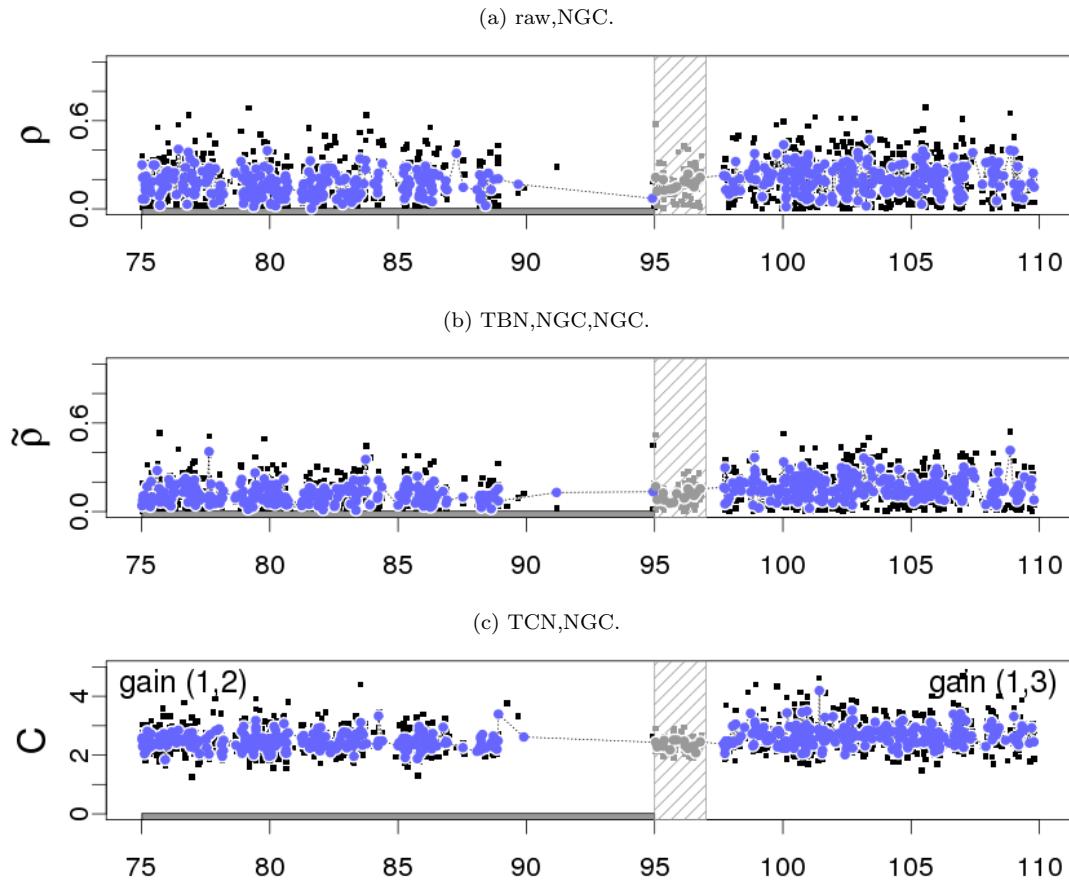


Figure 8: Decrease in Heterozygosity (DH) and total copy numbers for region TCGA-02-0001:Chr2@75-110,cp=96+/-1,s=1/4. Only heterozygous SNPs are plotted. There are 1023 loci of state gain (1,2) ("negatives") and 1023 loci of state gain (1,3) ("positives"), where the latter are highlighted with a solid bar beneath. In total 56 loci within the safety margin were excluded.

4.2 Allele B fraction density plots

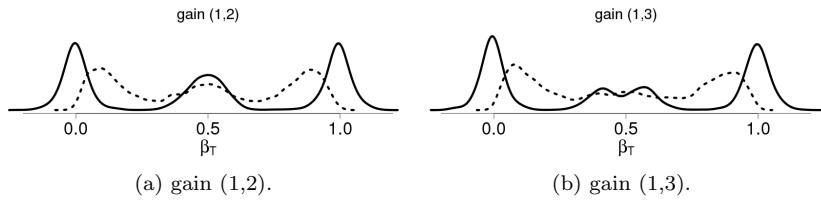


Figure 9: Density of raw (dashed lines) and TumorBoost-normalized (solid lines) allele B fractions for region TCGA-02-0001:Chr2@75-110,cp=96+/-1,s=1/4.

4.3 ROC curves

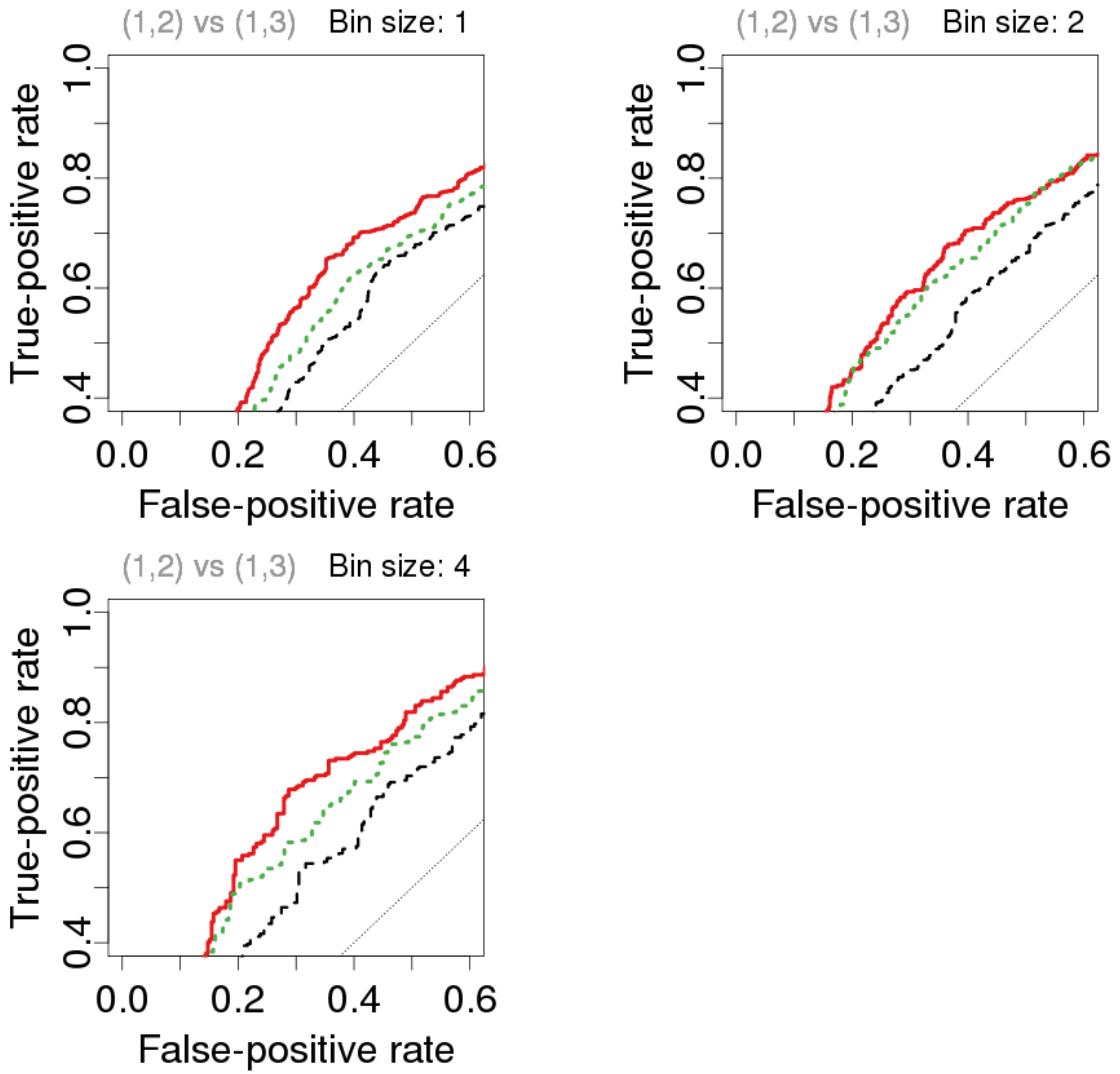


Figure 10: ROC curves for each preprocessing method at the full resolution as well as 2 different amounts of smoothing (using the `mean()` function) for region TCGA-02-0001:Chr2@75-110,cp=96+/-1,s=1/4. Legend: raw,NGC (dashed; #000000), TBN,NGC,NGC (solid; #E41A1C) and TCN,NGC (dotted; #4DAF4A).

4.4 (β_N, β_T) plots

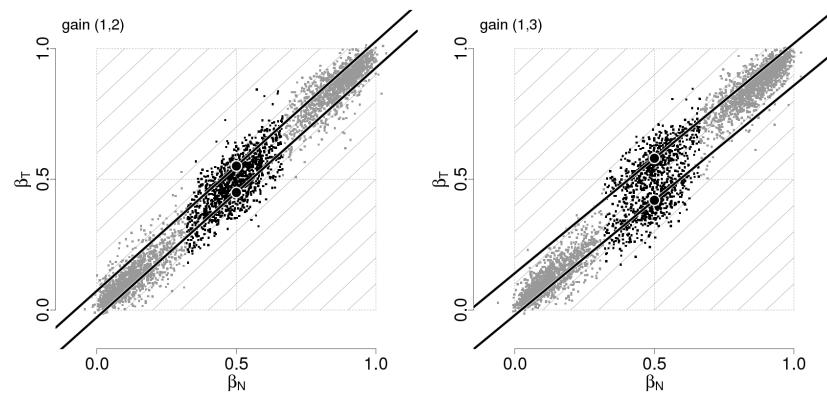


Figure 11: raw,NGC for region TCGA-02-0001:Chr2@75-110,cp=96+/-1,s=1/4.

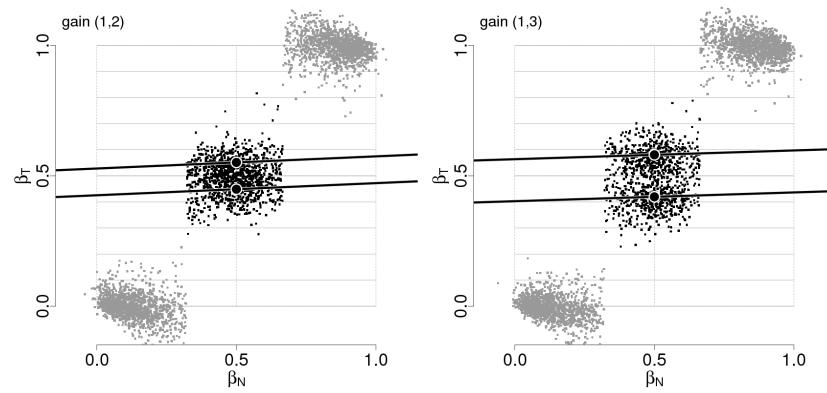


Figure 12: TBN,NGC,NGC for region TCGA-02-0001:Chr2@75-110,cp=96+/-1,s=1/4.

4.5 Allele-specific copy number estimates

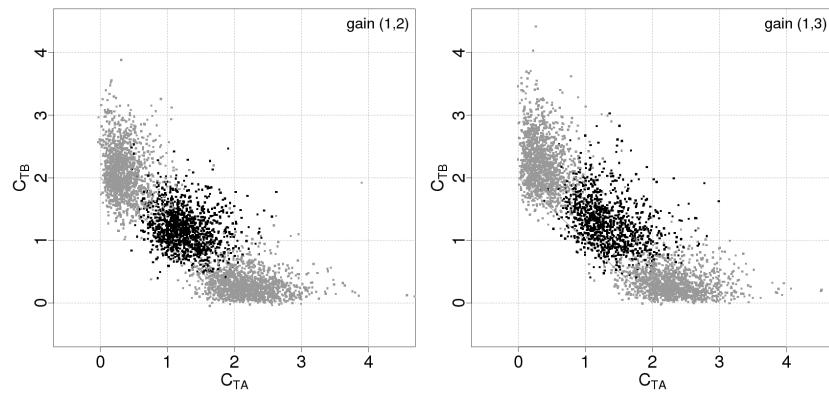


Figure 13: raw,NGC for region TCGA-02-0001:Chr2@75-110,cp=96+/-1,s=1/4.

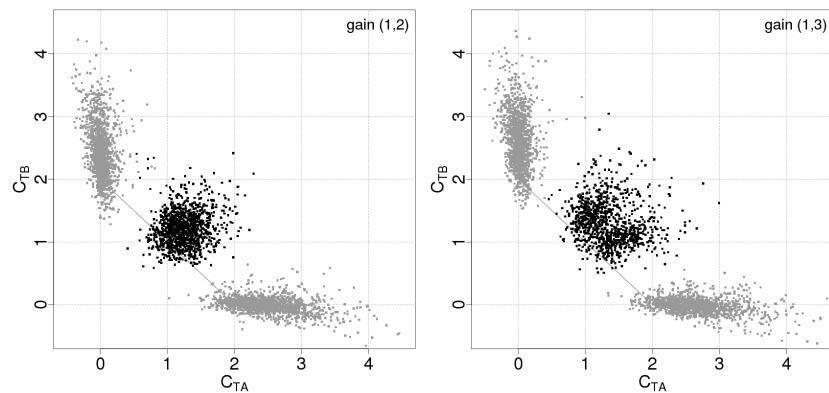


Figure 14: TBN,NGC,NGC for region TCGA-02-0001:Chr2@75-110,cp=96+/-1,s=1/4.

5 Region: TCGA-02-0001:Chr2@100-130,cp=110+/-1,s=4/0

5.1 Decrease in Heterozygosity (DH) and total copy-number tracks

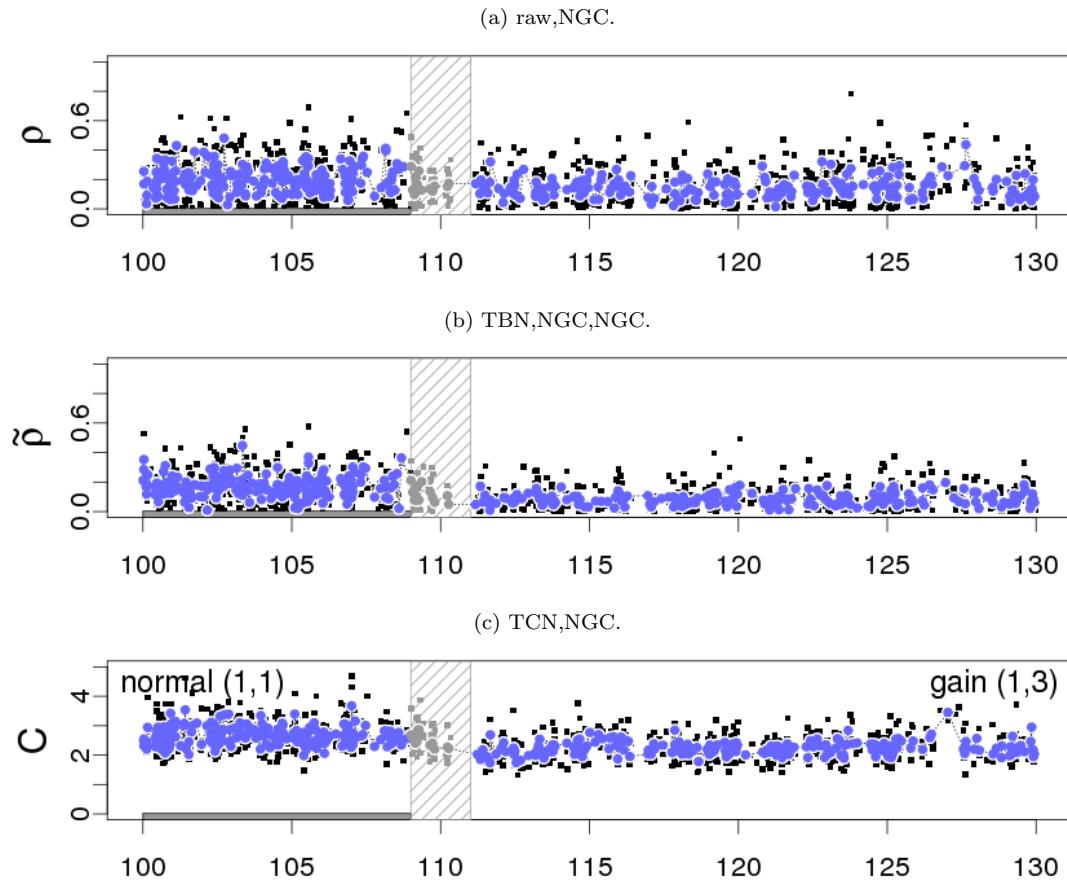


Figure 15: Decrease in Heterozygosity (DH) and total copy numbers for region TCGA-02-0001:Chr2@100-130,cp=110+/-1,s=4/0. Only heterozygous SNPs are plotted. There are 828 loci of state normal (1,1) ("negatives") and 828 loci of state gain (1,3) ("positives"), where the latter are highlighted with a solid bar beneath. In total 60 loci within the safety margin were excluded.

5.2 Allele B fraction density plots

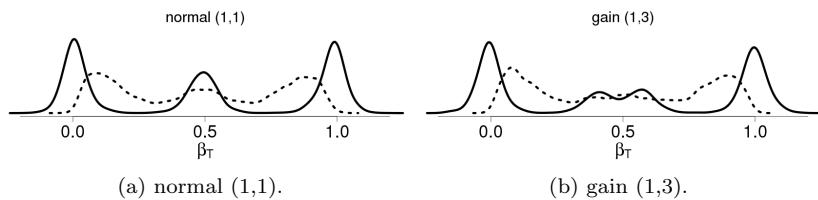


Figure 16: Density of raw (dashed lines) and TumorBoost-normalized (solid lines) allele B fractions for region TCGA-02-0001:Chr2@100-130,cp=110+/-1,s=4/0.

5.3 ROC curves

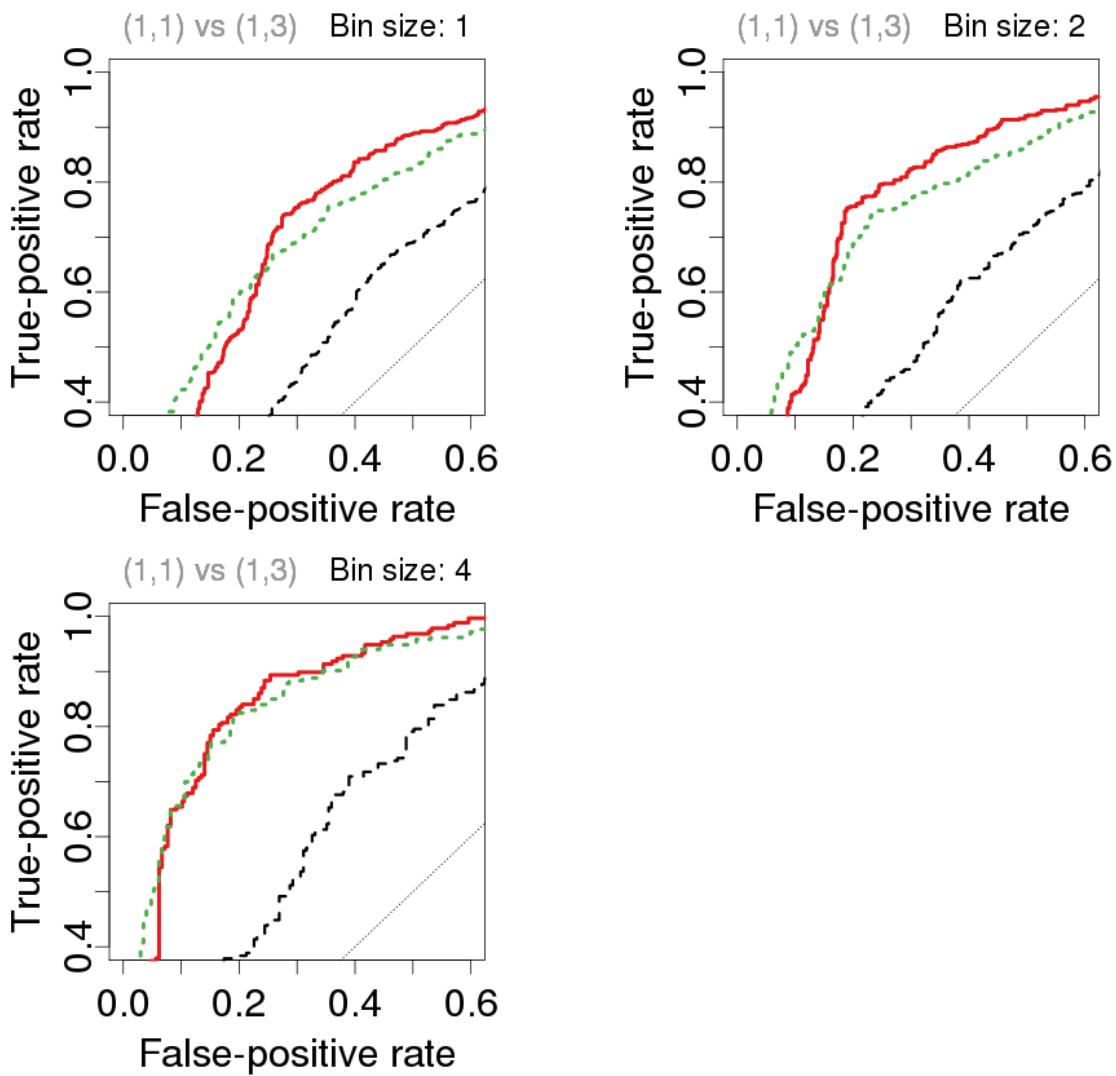


Figure 17: ROC curves for each preprocessing method at the full resolution as well as 2 different amounts of smoothing (using the `mean()` function) for region TCGA-02-0001:Chr2@100-130,cp=110+/-1,s=4/0. Legend: raw,NGC (dashed; #000000), TBN,NGC,NGC (solid; #E41A1C) and TCN,NGC (dotted; #4DAF4A).

5.4 (β_N, β_T) plots

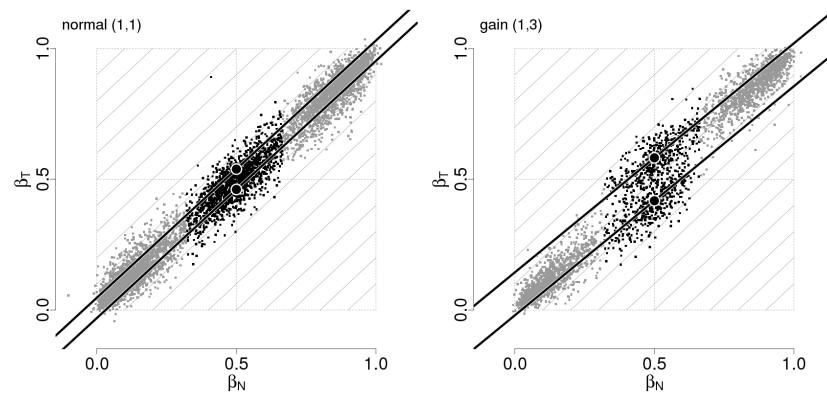


Figure 18: raw,NGC for region TCGA-02-0001:Chr2@100-130,cp=110+/-1,s=4/0.

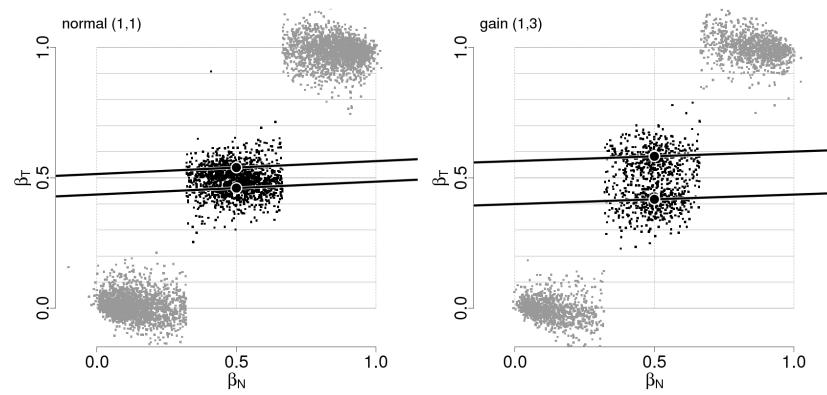


Figure 19: TBN,NGC,NGC for region TCGA-02-0001:Chr2@100-130,cp=110+/-1,s=4/0.

5.5 Allele-specific copy number estimates

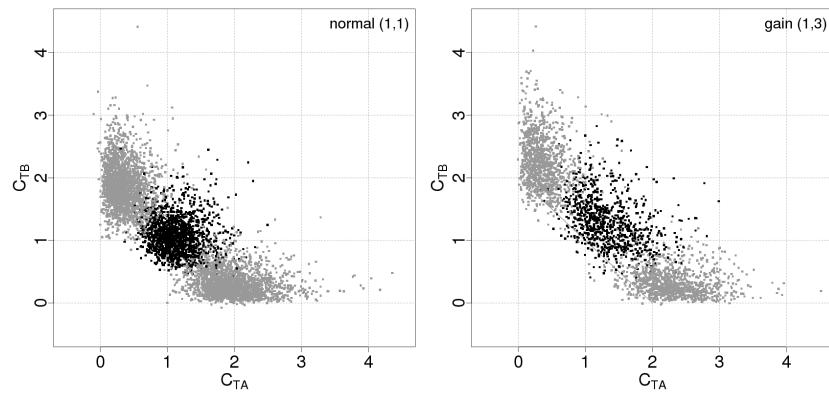


Figure 20: raw,NGC for region TCGA-02-0001:Chr2@100-130,cp=110+/-1,s=4/0.

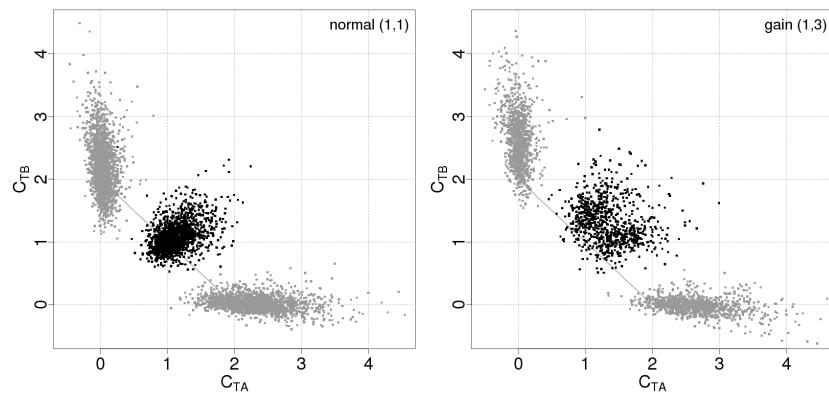


Figure 21: TBN,NGC,NGC for region TCGA-02-0001:Chr2@100-130,cp=110+/-1,s=4/0.

6 Region: TCGA-02-0001:Chr13@0-70,cp=45+/-1,s=0/2

6.1 Decrease in Heterozygosity (DH) and total copy-number tracks

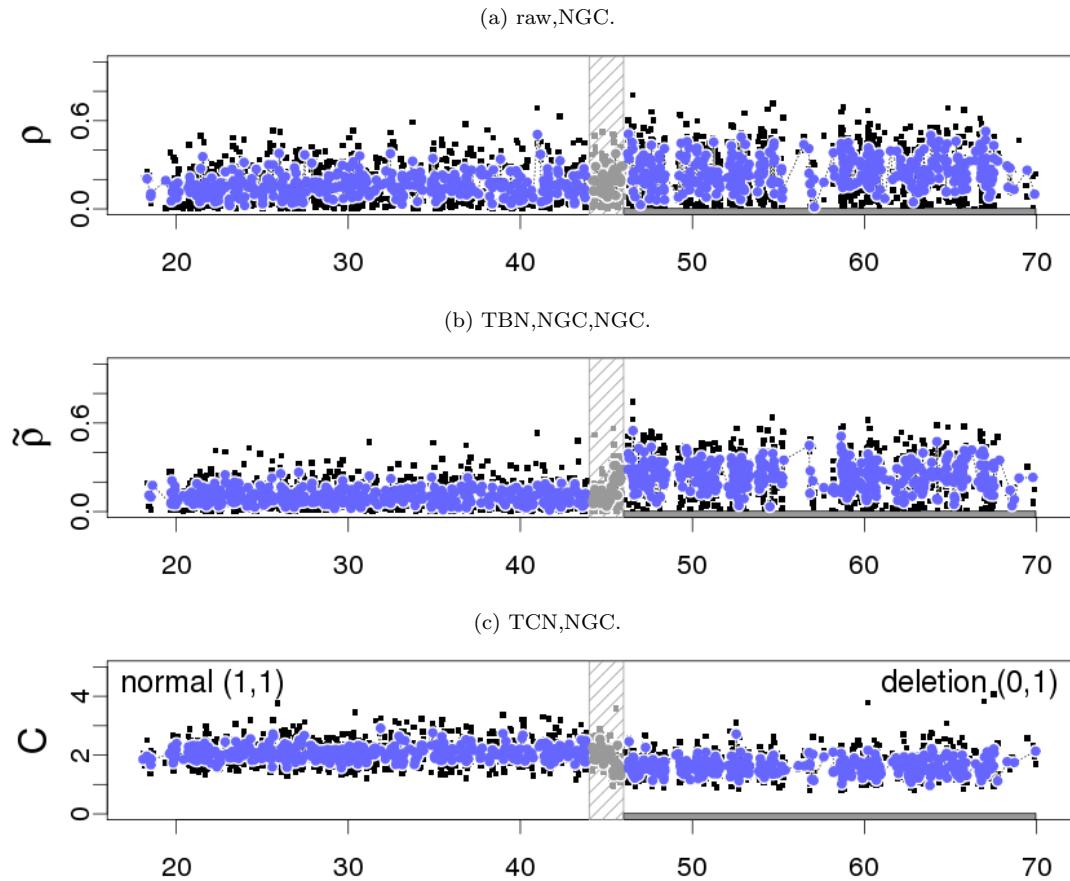


Figure 22: Decrease in Heterozygosity (DH) and total copy numbers for region TCGA-02-0001:Chr13@0-70,cp=45+/-1,s=0/2. Only heterozygous SNPs are plotted. There are 1566 loci of state normal (1,1) ("negatives") and 1566 loci of state deletion (0,1) ("positives"), where the latter are highlighted with a solid bar beneath. In total 147 loci within the safety margin were excluded.

6.2 Allele B fraction density plots

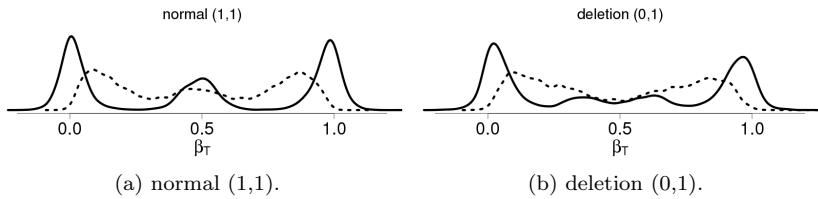


Figure 23: Density of raw (dashed lines) and TumorBoost-normalized (solid lines) allele B fractions for region TCGA-02-0001:Chr13@0-70,cp=45+/-1,s=0/2.

6.3 ROC curves

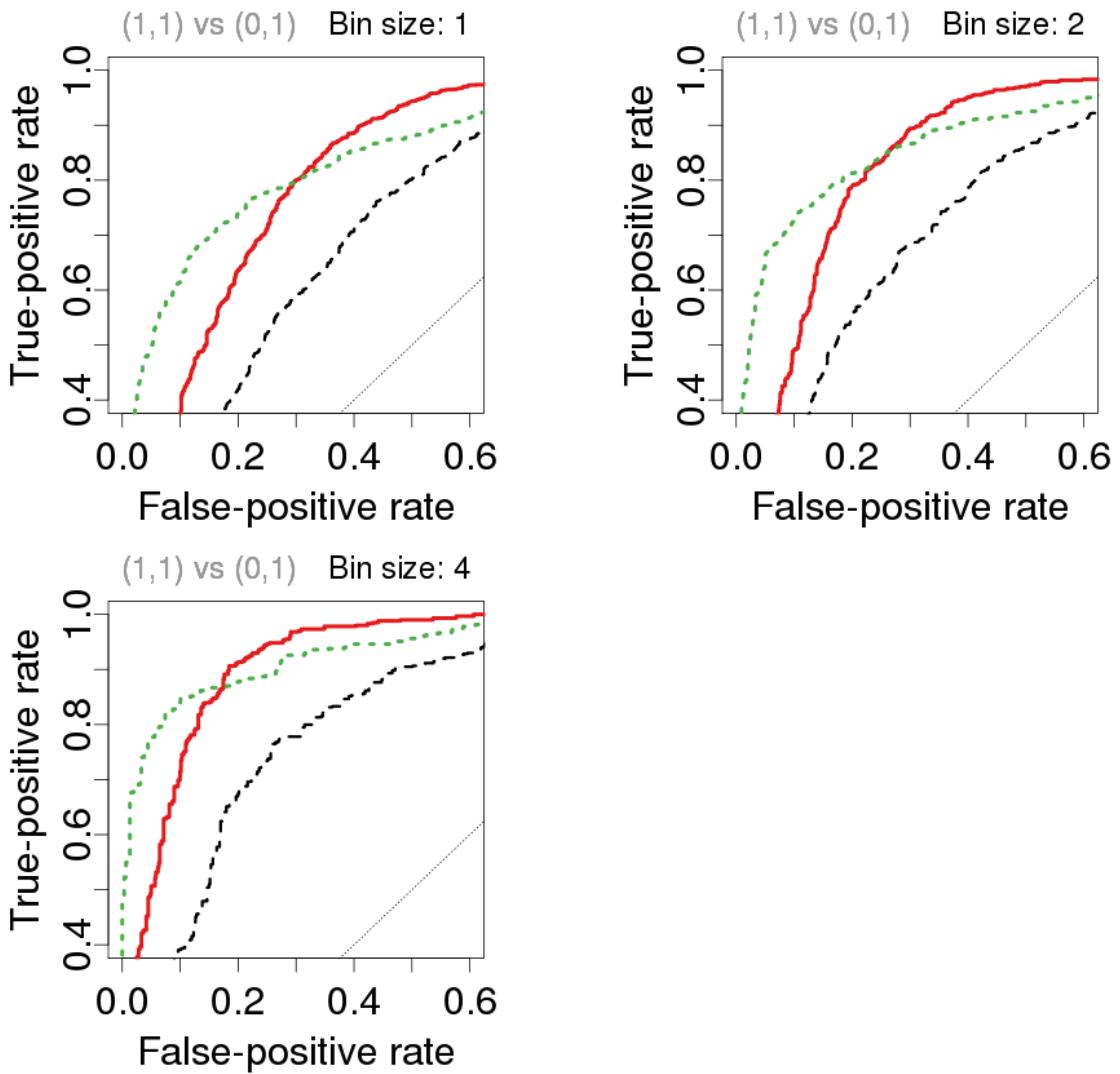


Figure 24: ROC curves for each preprocessing method at the full resolution as well as 2 different amounts of smoothing (using the `mean()` function) for region TCGA-02-0001:Chr13@0-70,cp=45+/-1,s=0/2. Legend: raw,NGC (dashed; `#000000`), TBN,NGC,NGC (solid; `#E41A1C`) and TCN,NGC (dotted; `#4DAF4A`).

6.4 (β_N, β_T) plots

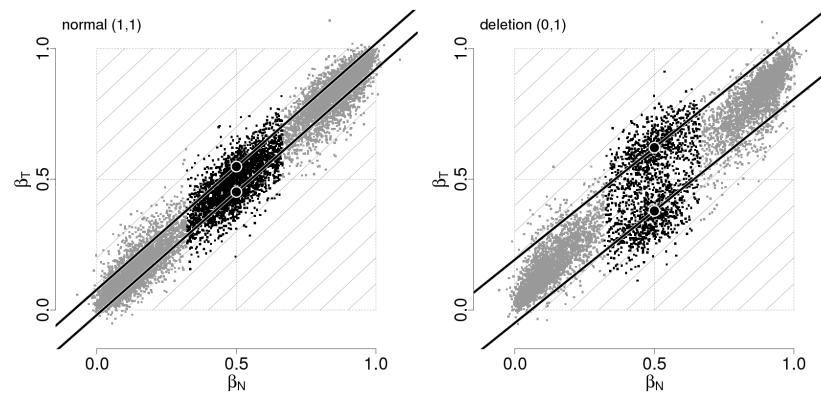


Figure 25: raw,NGC for region TCGA-02-0001:Chr13@0-70,cp=45+/-1,s=0/2.

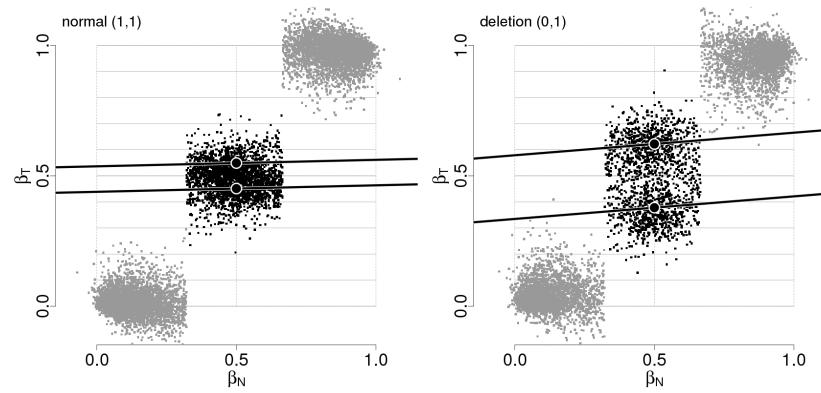


Figure 26: TBN,NGC,NGC for region TCGA-02-0001:Chr13@0-70,cp=45+/-1,s=0/2.

6.5 Allele-specific copy number estimates

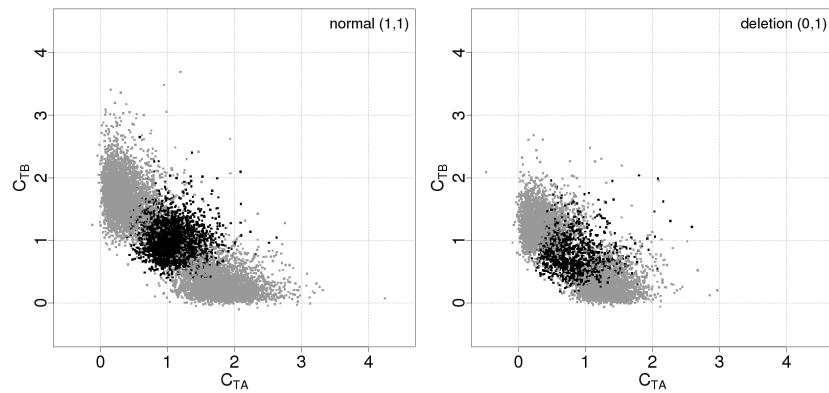


Figure 27: raw,NGC for region TCGA-02-0001:Chr13@0-70,cp=45+/-1,s=0/2.

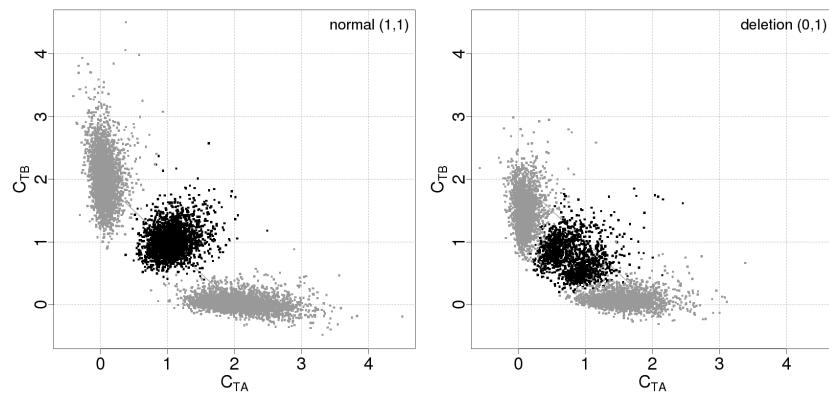


Figure 28: TBN,NGC,NGC for region TCGA-02-0001:Chr13@0-70,cp=45+/-1,s=0/2.

7 Bootstrap estimates of test statistics for all regions

| | 1/2 | 1/4 | 4/0 | 0/2 |
|-------------|--------------|-------------|--------------|--------------|
| raw,NGC | 11.352±1.053 | 2.958±0.965 | 4.582±1.048 | 8.635±1.124 |
| TBN,NGC,NGC | 17.171±1.257 | 6.965±1.015 | 10.960±1.245 | 14.306±1.098 |
| TCN,NGC | 20.770±1.241 | 5.200±0.971 | 11.095±1.056 | 15.041±1.275 |

Table 2: Student test statistics of the null hypothesis of equal mean before and after each PCN change point (heterozygous SNPs): raw or TumorBoost-normalized DH, and total copy number (last line). Mean \pm standard deviation across 100 samplings of 250 points (for each PCN state) from the original data set. The larger value, the more different the true means are.

References

- Bengtsson, H., Wirapati, P., and Speed, T. P. (2009). A single-array preprocessing method for estimating full-resolution raw copy numbers from all Affymetrix genotyping arrays including GenomeWideSNP 5 & 6. *Bioinformatics*, **25**(17), 2149–2156.

A Data files

A.1 Total copy numbers

```
$'TCGA,GBM,CRMAv2'  
AromaUnitTotalCnBinarySet:  
Name: TCGA  
Tags: GBM,CRMAv2  
Full name: TCGA,GBM,CRMAv2  
Number of files: 1  
Names: TCGA-02-0001 [1]  
Path (to the first file): rawCnData/TCGA,GBM,CRMAv2/GenomeWideSNP`6  
Total file size: 7.08 MB  
RAM: 0.00MB
```

A.2 Allele B fractions

```
$raw  
AromaUnitFracBCnBinarySet:  
Name: TCGA  
Tags: GBM,CRMAv2  
Full name: TCGA,GBM,CRMAv2  
Number of files: 1  
Names: TCGA-02-0001 [1]  
Path (to the first file): totalAndFracBData/TCGA,GBM,CRMAv2/GenomeWideSNP`6  
Total file size: 7.08 MB  
RAM: 0.00MB
```

```
$'TBN,NGC'  
AromaUnitFracBCnBinarySet:  
Name: TCGA  
Tags: GBM,CRMAv2,TBN,NGC  
Full name: TCGA,GBM,CRMAv2,TBN,NGC  
Number of files: 1  
Names: TCGA-02-0001 [1]  
Path (to the first file): totalAndFracBData/TCGA,GBM,CRMAv2,TBN,NGC/GenomeWideSNP`6  
Total file size: 7.08 MB  
RAM: 0.00MB
```

A.3 Genotype calls

```
$NGC  
AromaUnitGenotypeCallSet:  
Name: TCGA  
Tags: GBM,CRMAv2,NGC  
Full name: TCGA,GBM,CRMAv2,NGC  
Number of files: 1  
Names: TCGA-02-0001 [1]  
Path (to the first file): callData/TCGA,GBM,CRMAv2,NGC/GenomeWideSNP`6  
Total file size: 3.54 MB  
RAM: 0.00MB
```

B Session information

This report was automatically generated using the R.rsp package.

- R version 2.10.1 (2009-12-14), x86_64-unknown-linux-gnu
- Locale: LC_CTYPE=en_US.UTF-8, LC_NUMERIC=C, LC_TIME=en_US.UTF-8, LC_COLLATE=en_US.UTF-8, LC_MONETARY=C, LC_MESSAGES=en_US.UTF-8, LC_PAPER=en_US.UTF-8, LC_NAME=C, LC_ADDRESS=C, LC_TELEPHONE=C, LC_MEASUREMENT=en_US.UTF-8, LC_IDENTIFICATION=C
- Base packages: base, datasets, graphics, grDevices, methods, stats, utils
- Other packages: aroma.cn.eval 0.1.1, aroma.core 1.5.0, aroma.light 1.15.1, digest 0.4.2, MASS 7.3-4, matrixStats 0.1.9, R.cache 0.2.0, RColorBrewer 1.0-2, R.filesets 0.8.0, R.menu 0.1.0, R.methodsS3 1.1.0, R.oo 1.6.7, R.rsp 0.3.6, R.utils 1.3.3, xtable 1.5-6
- Loaded via a namespace (and not attached): affxparser 1.18.0, tools 2.10.1